

# ExaGPT: Example-Based Machine-Generated Text Detection for Human Interpretability

Anonymous ACL submission

## Abstract

001 Detecting texts generated by Large Language  
002 Models (LLMs) could cause grave mistakes  
003 due to incorrect decisions, such as undermin-  
004 ing student’s academic dignity. LLM text de-  
005 tection thus needs to ensure the interpretability  
006 of the decision, which can help users judge  
007 how reliably correct its prediction is. When  
008 humans verify whether a text is human-written  
009 or LLM-generated, they intuitively investigate  
010 with which of them it shares more similar spans.  
011 However, existing interpretable detectors are  
012 not aligned with the human decision-making  
013 process and fail to offer evidence that users  
014 easily understand. To bridge this gap, we intro-  
015 duce **ExaGPT**, an interpretable detection ap-  
016 proach grounded in the human decision-making  
017 process for verifying the origin of a text. Exa-  
018 GPT identifies a text by checking whether it  
019 shares more similar spans with human-written  
020 vs. with LLM-generated texts from a datastore.  
021 This approach can provide similar span exam-  
022 ples that contribute to the decision for each  
023 span in the text as evidence. Our human evalua-  
024 tion demonstrates that providing similar span  
025 examples contributes more effectively to judg-  
026 ing the correctness of the decision than exist-  
027 ing interpretable methods. Moreover, extensive ex-  
028 periments in four domains and three generators  
029 show that ExaGPT massively outperforms prior  
030 interpretable detectors by up to +37.0 points of  
031 accuracy at a false positive rate of 1%. We will  
032 release our code after acceptance.

## 033 1 Introduction

034 LLMs can yield human-like texts in response to var-  
035 ious textual instructions (OpenAI, 2023b; Touvron  
036 et al., 2023). Ironically, the powerful generative  
037 capability has resulted in various misuses of LLMs,  
038 such as cheating in student homework assignments  
039 and mass-producing fake news (Tang et al., 2023;  
040 Wu et al., 2023). Such abuse of LLMs has sparked  
041 the demand for discerning LLM-generated texts  
042 from human-written ones.

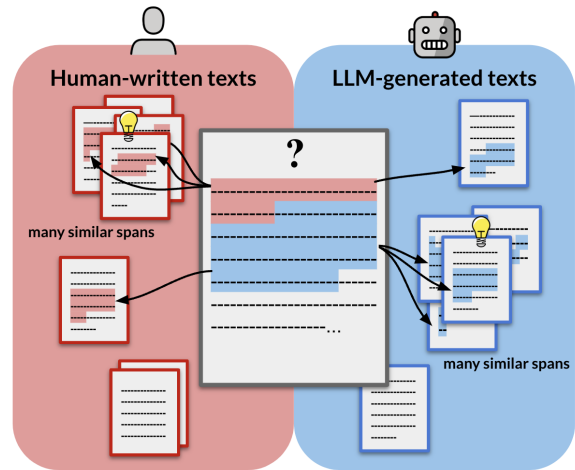


Figure 1: Identifying the author of a text (human vs. LLM) by examining if it shares more similar spans, including verbatim overlaps and semantically similar spans, with human-written vs. LLM-generated texts.

043 While recent powerful detectors (Mitchell et al.,  
044 2023; Su et al., 2023a; Koike et al., 2024; Hans  
045 et al., 2024; Verma et al., 2024) can help prevent  
046 potential misuse of LLMs, misclassifications could  
047 lead to severe consequences. For instance, web  
048 content writers have recently been at risk of los-  
049 ing their careers because of false-positive classi-  
050 fication (Gizmodo, 2024). In school education,  
051 incorrect detection results might ruin students’ aca-  
052 demic dignity (OpenAI, 2023a; Bloomberg, 2024).  
053 At the same time, it is extremely difficult, if not im-  
054 possible, to develop a perfect detector with 100%  
055 accuracy in such real-world scenarios. There re-  
056 main edge cases where human-written texts can  
057 be misidentified as LLM-generated and vice versa.  
058 Thus, it is crucial to create a detector that provides  
059 interpretable evidence, allowing users to judge how  
060 reliably correct the detection results are and recog-  
061 nize potential misclassifications (Tang et al., 2023;  
062 Ji et al., 2024).

063 Most detectors lack the interpretability of their

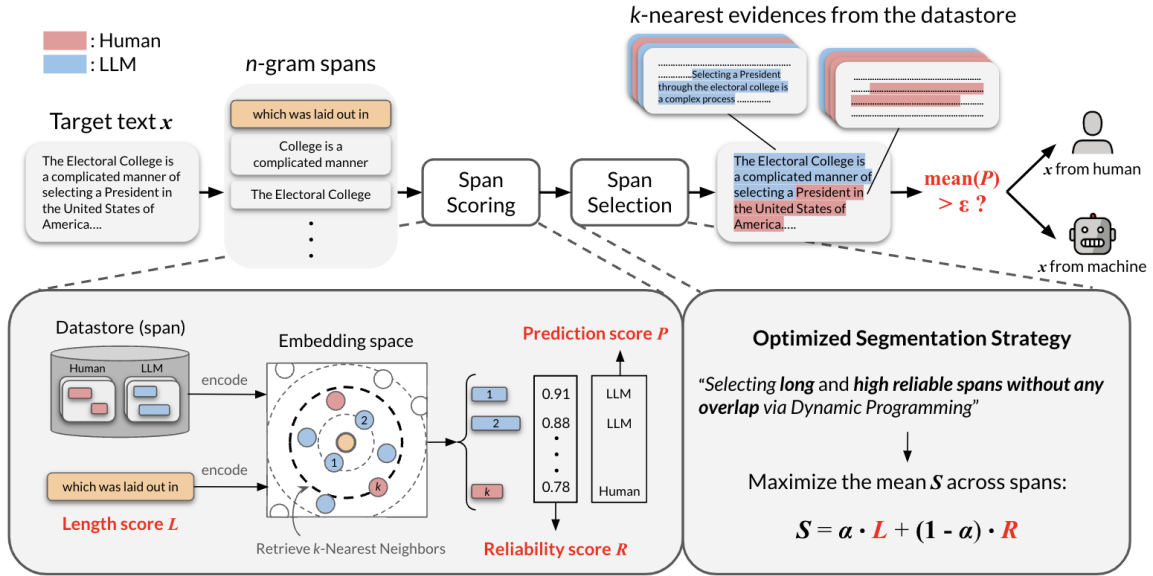


Figure 2: Overview of ExaGPT. It detects the author of a text by examining whether the text shares more similar spans with human-written texts vs. with LLM-generated texts from a datastore.

064 decisions, outputting only binary labels of who au- 095  
 065 thored the text. There are few studies on the inter- 096  
 066 pretability of the detection. Gehrmann et al. (2019) 097  
 067 color-highlighted the tokens with high probability 098  
 068 under the predicted distribution of LMs. Mitrović 099  
 069 et al. (2023); Wang et al. (2024) showed which part 100  
 070 of a text contributed to a decision based on predic- 101  
 071 tion shifts via perturbations to the text. Yang et al. 102  
 072 (2023) provided the  $n$ -gram overlaps between the 103  
 073 original text and re-prompted ones generated by 104  
 074 LLMs. Here, humans intuitively judge whether a 105  
 075 text is human-written or LLM-generated by assess- 106  
 076 ing with which source it shares more *similar spans*, 107  
 077 including verbatim overlaps and semantically simi- 108  
 078 lar spans (Maurer et al., 2006; Barrón-Cedeño et al., 109  
 079 2013). However, current detectors are not aligned 110  
 080 with the human decision-making process (Figure 1) 111  
 081 and fail to yield sufficiently interpretable evidence 112  
 082 for users. 113

083 Motivated by this gap, we present **ExaGPT**, an 114  
 084 interpretable detection method based on the human 115  
 085 decision-making process of verifying the origin of 116  
 086 a text. In particular, ExaGPT makes a prediction 117  
 087 by examining whether the text shares more similar 118  
 088 spans with human-written vs. with LLM-generated 119  
 089 texts from a datastore. This approach can provide 120  
 090 similar span examples that contribute to the deci- 121  
 091 sion for each span in the text as interpretable 122  
 092 evidence. To present interpretable span-segmented 123  
 093 text as a final result, we apply a dynamic program- 124  
 094 ming algorithm and determine the optimal span

break. It balances the long span length and its high 095  
 frequency with the datastore (i.e., many similar 096  
 phrases to the span exist in the datastore). The 097  
 similarity of the retrieved spans to each span in the 098  
 target text can help users judge the reliability of the 099  
 detection result. 100

To evaluate the interpretability of LLM detec- 101  
 tion, we conducted a human evaluation of how well 102  
 people can infer the correctness of the detection 103  
 from the detector’s evidence, and we found that 104  
 providing similar span examples contributes more 105  
 effectively to judging the correctness of the detec- 106  
 tion than existing interpretable methods. Moreover, 107  
 extensive experiments in four domains and three 108  
 generators showed that ExaGPT massively outper- 109  
 forms prior interpretable and powerful detectors 110  
 by up to +37.0 points accuracy, even at a constant 111  
 false positive rate of 1%. From these results, we 112  
 observe that ExaGPT achieves high interpretability 113  
 in its detection result and also high detection 114  
 performance. 115

## 2 Methodology 116

ExaGPT classifies a text based on whether it shares 117  
 more similar spans with human-written or with 118  
 LLM-generated texts from a datastore. As a final 119  
 result, ExaGPT offers the span-segmented text 120  
 where each span is accompanied by similar span 121  
 examples that contribute to the decision. Figure 2 122  
 illustrates the workflow of ExaGPT, which has two 123  
 phases: **Span Scoring** and **Span Selection**. In 124

the first phase, we mainly investigate whether each span in the target text shares more similar spans with human-written or LLM-generated texts from a datastore. Meanwhile, we calculate scores for each span, which we use in the second phase (§2.1). In the second phase, we primarily decide the optimal span segmentation to aid users’ understanding of the final result. Specifically, we apply a dynamic programming (DP) algorithm with the scores from the first phase to find the span boundaries, balancing span length and its frequency within the datastore (§2.2). Finally, we detect the target text based on the selected spans and we provide similar span examples for each target span as evidence (§2.3). We will go into further details below.

## 2.1 Span Scoring with $k$ -NN Search

Given a target text  $x$  to be classified, we define an  $n$ -gram span in the text  $x$  as  $x_{i:i+n}$ , which is any continuous sequence of  $n$  tokens starting in the  $i$ -th token. For each  $n$ -gram target span  $x_{i:i+n}$ , we retrieve the top- $k$  most similar<sup>1</sup>  $n$ -gram spans  $s_j$  ( $j \in \{1, \dots, k\}$ ) from the datastore, with each original label and similarity  $\{(s_j, l_j, c_j)\}_{j=1}^k$ . Here,  $l_j$  is Human when the span  $s_j$  is part of a human-written text, or LLM when the span  $s_j$  is a part of a LLM-generated text.  $c_j$  is the similarity between the target span  $x_{i:i+n}$  and each retrieved span  $s_j$ .

Consequently, we calculate the following metrics for each target span  $x_{i:i+n}$ : *length score*  $L$ , *reliability score*  $R$ , and *prediction score*  $P$ . The length score  $L$  is the number of tokens in the target span:

$$L(x_{i:i+n}) = n \quad (1)$$

The reliability score  $R$  is the mean similarity  $c_j$  between the target span and each retrieved span:

$$R(x_{i:i+n}) = \frac{\sum_{j=1}^k c_j}{k} \quad (2)$$

The reliability score  $R$  indicates how many similar spans exist in the datastore for the target span. The prediction score  $P$  is a ratio of LLM label in the original labels  $l_j$  of the retrieved spans:

$$P(x_{i:i+n}) = \frac{\sum_{j=1}^k \mathbb{1}(l_j = \text{LLM})}{k}. \quad (3)$$

The prediction score  $P$  indicates whether the target span shares more similar spans with human-written vs. with LLM-generated texts in the datastore.

<sup>1</sup>We encode the target span, and all spans in the datastore into the same embedding space. We then perform  $k$ -nearest neighbor ( $k$ -NN) search based on the cosine similarity of each two span embeddings. See more details in §3.1.

---

## Algorithm 1 Span Segmentation Optimization

---

**Input:** Target text  $x$ ; Length of target text  $m$ ; Length score  $L$ ; Reliability score  $R$ ; Maximum length of  $n$ -gram span  $N$ ; Hyper-parameter  $\alpha$

**Output:** List of selected  $n$ -grams  $T$   
 $\text{dp}[0, \dots, m-1] \leftarrow [([0], \text{None})] * m$

```

for  $i = 1$  to  $m$  do
  for  $j = \min(i - N, 0)$  to  $i$  do
     $l, r \leftarrow L^{\text{std}}(x_{j:i}), R^{\text{std}}(x_{j:i})$ 
     $\text{scores} \leftarrow \text{dp}[j][0] + [\alpha l + (1 - \alpha)r]$ 
     $s_{\text{cand}} \leftarrow \text{average}(\text{scores})$ 
    if  $\text{average}(\text{dp}[i][0]) < s_{\text{cand}}$  then
       $\text{dp}[i] \leftarrow (\text{scores}, j)$ 
    end if
  end for
end for

```

Traverse dp backward and collect span breaks  
**return** List of selected  $n$ -grams  $T$

---

## 2.2 Span Selection with a DP Algorithm

In this phase, we select spans  $T = [t_1, \dots, t_H]$  in the target text  $x$ , so that the text is segmented without overlaps as a final result:

$$x = t_1 \oplus t_2 \oplus \dots \oplus t_H, \quad (4)$$

$$t_i \cap t_j = \emptyset \quad (i, j \in \{1, \dots, H\}, i \neq j)$$

To facilitate users’ understanding of the final result, we optimize the span segmentation that includes longer and more similar spans with ones from the datastore. Algorithm 1 describes our dynamic programming strategy to find the best span break. Formally, we select spans  $T$  to maximize the score  $S$  across the spans in the target text:

$$S(T) = \frac{\sum_{h=1}^H \{\alpha L^{\text{std}}(t_h) + (1 - \alpha)R^{\text{std}}(t_h)\}}{H}. \quad (5)$$

Here,  $L^{\text{std}}(t_h)$  and  $R^{\text{std}}(t_h)$  are the normalized<sup>2</sup> versions of the length score  $L$  and the reliability score  $R$  of the span  $t_h$ , respectively.  $\alpha$  is an interpolation coefficient ranging from 0.0 to 1.0.  $\alpha$  determines the relative contribution of the length score and the reliability score to the span segmentation.

## 2.3 Overall Detection with Evidence

Given a sequence of the selected spans  $T$  each with a prediction score for the target text  $x$ , ExaGPT

<sup>2</sup>To align the scales of the length score and the reliability score, each score is normalized using the mean and the variance in the validation split of our dataset.

identifies a text based on the mean prediction score:

$$P_{\text{overall}} = \frac{\sum_{h=1}^H P(t_h)}{H}. \quad (6)$$

ExaGPT classifies a text as LLM if  $P_{\text{overall}}$  exceeds a detection threshold  $\epsilon$ , and otherwise as Human. As evidence of the decision, ExaGPT provides retrieved top- $k$  similar spans for each span in the text:

$$E = [(t_h, [s_h^1, \dots, s_h^k])]_{h=1}^H. \quad (7)$$

The similarity of the retrieved spans to each span in the target text can help users judge how reliably correct the detection result is.

### 3 Experiments and Results

#### 3.1 Overall Setup

**Metrics.** To assess the detection performance, we use the AUROC score, which is widely used in studies on LLM detection. However, it is only useful to observe the overall behavior of a detector through all possible thresholds. In practice, it is critical to minimize the false positive classification, i.e., wrongly identifying human-written texts as LLM-generated. We thus report the detection accuracy with a threshold by fixing the false-positive rate (FPR) at 1%, which is an evaluation stream among recent robustness studies (Krishna et al., 2023; Hans et al., 2024; Dugan et al., 2024).

**Datasets.** We use the M4 dataset (Wang et al., 2024), which is a large-scale detection benchmark consisting of pairs of human-written and LLM-generated texts across multiple languages, domains, and generators. Our experiments use the English subset, including 3,000 pairs of human-written and LLM-generated texts from each combination of four domains: Wikipedia, Reddit, WikiHow, and arXiv, and three generators: ChatGPT, GPT-4 as closed-source LLMs, and Dolly-v2 (Conover et al., 2023) as open-source LLMs. For each combination, we split the dataset into three parts: train/valid/test with 2,000/500/500 pairs, respectively.

**Baselines.** We compare ExaGPT to three strong and interpretable detectors (as detailed in §5): RoBERTa with SHAP (Mitrović et al., 2023), LR-GLTR (Wang et al., 2024), and DNA-GPT (Yang et al., 2023). The first one is a supervised classifier based on RoBERTa (Liu, 2019), which we fine-tune for detection on our train split. Similarly, we train the LR-GLTR detector on our train

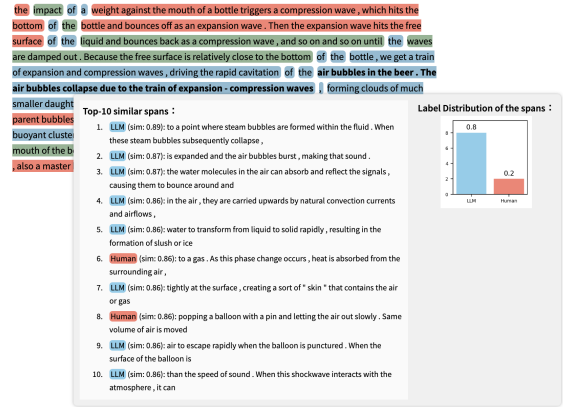


Figure 3: User interface of ExaGPT. Hovering over a text span displays the tooltip about the retrieved similar spans each with the similarity to the span and the original label distribution.

split with selected and hand-crafted GLTR features (Gehrmann et al., 2019), following (Wang et al., 2024). The hyper-parameter settings for training both RoBERTa and LR-GLTR are aligned with (Wang et al., 2024). Further configurations of the baselines are in Appendix A.

**ExaGPT.** In the span scoring phase, ExaGPT leverages our train split as the datastore for each combination of domains and generators. We consider the size of  $n$ -gram to be from 1 to 20 throughout the entire dataset. We embed the target span and all spans in the datastore into the same vector space using BERT<sup>3</sup>. For a span embedding, we feed a text into BERT and take the mean second-layer<sup>4</sup> hidden outputs of tokens included in the span. We retrieve the top- $k$  ( $=10$ )<sup>5</sup> most similar spans from the datastore for each target span via  $k$ -NN search using the FAISS (Johnson et al., 2017).

In the span selection phase, we select the optimal  $\alpha$  from values between 0.0 and 1.0 at 0.125 intervals, where ExaGPT exhibits the best detection performance in our validation split. The  $\alpha$  is constant through our evaluation of the interpretability and the detection performance of ExaGPT.

<sup>3</sup><https://huggingface.co/google-bert/bert-large-uncased>

<sup>4</sup>We select the layer where the  $k$ -NN spans are similar to the target span well-balanced lexically and semantically, enhancing its interpretability in our pilot study.

<sup>5</sup>We choose the value of  $k$  so that ExaGPT shows favorable detection performance over smaller values in our pilot study and does not reduce the interpretability of its evidence. Since ExaGPT presents retrieved spans as evidence, keeping  $k$  small helps users assess detection correctness based on a manageable amount of information.

**Human Evaluation on Interpretability.** We assess the interpretability of detectors through human evaluation, since it is crucial for a detector to provide evidence that allows users to judge how reliable the detection result is. Therefore, we first design a human evaluation that tests whether the provided evidence *actually helps* users judge whether the detection is correct, which is a practical aspect overlooked in prior detection work. Specifically, participants are shown the detection evidence and asked to judge whether the detection is correct. Consequently, our interpretability metric is defined as the accuracy of human judgments of detection correctness based on the evidence.

For each detector, we evaluate 96 samples<sup>6</sup> from our test split in all combinations of domains and generators so that the ratio of correct and incorrect detections<sup>7</sup> is even. In our human evaluation, four annotators, including one MSc student, one PhD student, and two researchers working in natural language processing, were provided with different samples. Figure 3 shows the user interface of ExaGPT in our human evaluation. The spans are highlighted<sup>8</sup> in red, green, and blue for which prediction score  $P$  is lower than 0.5 (human-written), equal to 0.5 (neither), and higher than 0.5 (LLM-generated), respectively. The participants identify the correctness of the detection by mainly investigating similar span examples for each span in the text. We elaborate on the detection evidence of each baseline detector in Appendix B.

## 3.2 Results

**Detection Interpretability.** Table 1 presents the difference in the accuracy of human judgments on the detection correctness based on evidence across baseline detectors and ExaGPT. The accuracy of human judgments on ExaGPT is relatively higher compared to baseline detectors by up to +13.6 points. This indicates that ExaGPT offers more interpretable evidence than other baselines, helping humans judge the correctness of detections more effectively. Here, DNA-GPT also offers  $n$ -gram span overlaps between the target text and the re-generated LLM texts from the truncated part as

<sup>6</sup>The 96 samples for each detector consist of two samples (one correct and one incorrect) across four domains and three generators, distributed among four participants.

<sup>7</sup>We focus on the setting of the 1% FPR threshold based on practical scenarios.

<sup>8</sup>ExaGPT performs the overall detection rather than detecting each span individually. However, for better readability, each span is color-highlighted on its prediction score.

Detector	ACC. of Human Judgements (%) $\uparrow$
RoBERTa	47.9
LR-GLTR	57.3
DNA-GPT	53.1
ExaGPT	<b>61.5</b>

Table 1: Comparison of the accuracy (ACC.) of human judgments on the correctness of detections based on evidence across baseline detectors and ExaGPT. Higher accuracy implies that the detector provides more interpretable evidence to users.

evidence. The comparison of the human evaluation score between DNA-GPT and ExaGPT suggests that providing not only simple overlaps but also semantically similar spans contributes to better interpretability. We further investigate how the similarity between the target span and retrieved spans correlates with the correctness of the detection of ExaGPT in §4.

**Detection Performance.** Table 2 presents the differences in detection performance between baseline detectors and ExaGPT across domains and generators. The detection performance includes AUROC and the accuracy at 1% FPR. Overall, ExaGPT consistently demonstrates detection performance on par with or better than baselines, including supervised classifiers. Specifically, on accuracy at 1% FPR, ExaGPT achieves the best average detection performance on all three generators, outperforming baselines by a large margin of up to +37.0 points. This suggests that ExaGPT is the most effective detector in practical scenarios, where we need to minimize the false positives.

In summary, ExaGPT achieved both superior interpretability of the detection and exceptional detection performance compared to previous interpretable detectors.

## 4 Analysis

**What Makes ExaGPT Interpretable.** Our human evaluations demonstrate that ExaGPT provides highly interpretable evidence compared to prior detectors. To explore the reason for this, we investigate the difference in the characteristics of the selected spans as a final output between correct and incorrect predictions by ExaGPT. We focus on span length and mean similarity between each target span and the retrieved spans (reliability score  $R$ ), which are prioritized in the span selection. We randomly select 1,000 correct and 1,000 incorrect ExaGPT predictions on our test splits across all

Generator	Detector	Wikipedia		Reddit		WikiHow		arXiv		Average	
		AUROC	ACC.	AUROC	ACC.	AUROC	ACC.	AUROC	ACC.	AUROC	ACC.
ChatGPT	RoBERTa	<b>100.0</b>	77.1	<b>99.8</b>	61.0	<b>100.0</b>	50.0	<b>100.0</b>	87.3	<b>100.0</b>	68.9
	LR-GLTR	95.0	60.0	<u>99.4</u>	<b>94.0</b>	97.5	85.8	<u>99.8</u>	<b>97.7</b>	97.9	<u>84.4</u>
	DNA-GPT	84.8	49.4	92.3	62.9	99.4	93.5	89.0	59.9	91.4	66.4
	ExaGPT	<u>98.6</u>	<b>92.3</b>	98.9	<u>86.6</u>	<u>99.5</u>	<b>96.0</b>	99.6	<u>95.8</u>	<u>99.2</u>	<b>92.7</b>
GPT-4	RoBERTa	<b>100.0</b>	<b>87.8</b>	<b>100.0</b>	66.4	<b>100.0</b>	77.4	<b>100.0</b>	68.6	<b>100.0</b>	75.1
	LR-GLTR	97.8	85.7	<u>99.6</u>	<b>97.2</b>	94.8	77.8	<b>100.0</b>	98.5	98.1	89.8
	DNA-GPT	40.3	48.1	<u>71.9</u>	68.6	44.6	49.9	72.2	54.4	57.3	<u>55.3</u>
	ExaGPT	<u>98.3</u>	<u>87.3</u>	99.3	<u>91.1</u>	<u>98.8</u>	<b>92.2</b>	<u>99.7</u>	<b>98.7</b>	<u>99.0</u>	<b>92.3</b>
Dolly-v2	RoBERTa	<b>100.0</b>	61.8	<b>100.0</b>	50.0	<b>100.0</b>	70.8	<b>100.0</b>	<b>82.8</b>	<b>100.0</b>	66.4
	LR-GLTR	79.7	57.7	95.3	<b>79.0</b>	72.4	55.0	<u>93.7</u>	<u>78.2</u>	85.3	<u>67.5</u>
	DNA-GPT	68.0	61.5	67.5	66.1	87.7	<b>82.3</b>	<u>64.9</u>	57.7	72.0	66.9
	ExaGPT	<u>85.8</u>	<b>63.8</b>	<u>96.2</u>	<u>76.6</u>	<u>94.3</u>	<u>75.6</u>	85.2	67.3	<u>90.4</u>	<b>70.8</b>

Table 2: **Comparison of detection performances** of ExaGPT and baseline detectors on texts from various domains and generators. *ACC.* indicates the detection accuracy at 1% FPR. **Bold** and Underline indicate the best and runner-up performance for each combination of domains and generators.

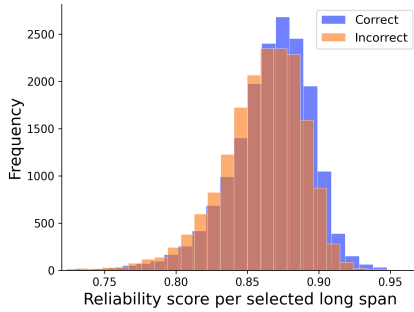


Figure 4: Reliability score distributions of long spans ( $n \geq 10$ ) in correct and incorrect samples of ExaGPT.

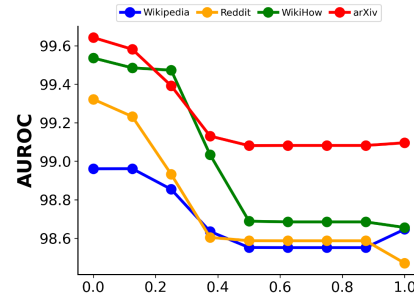


Figure 5: Impact of  $\alpha$  on the detection performance of ExaGPT, using ChatGPT as a generator.

combinations of domains and generators.

Figure 4 presents the reliability score distributions of long spans ( $n \geq 10$ ) for correct and incorrect samples. A rightward shift indicates that correct samples of ExaGPT include more long spans with higher reliability scores than incorrect ones. This suggests that offering long spans with high reliability scores helps users judge the correctness of the detections. Table 7 presents examples of long spans ( $n = 19$ ) with high reliability scores for a target span retrieved by ExaGPT. We can see that the retrieved spans are well-balanced between lexical and semantic similarity to the target span. Due to space limitations, the full table is provided in Appendix C.

**Impact of  $\alpha$ .** In our experiments, we determine the optimal coefficient  $\alpha$  for ExaGPT (as used in Eq. 5) based on the best detection performance on the validation split. To examine the robustness of ExaGPT to the choice of  $\alpha$ , we analyze how detection performance varies as  $\alpha$  changes.

Figure 5 depicts the relationship between  $\alpha$  and

the detection performance of ExaGPT, evaluated on ChatGPT-generated text across four domains.  $\alpha$  ranges from 0.0 to 1.0 in increments of 0.125. We observe that larger values of  $\alpha$  generally lead to lower detection performance. This suggests that placing greater weight on the reliability score (i.e., selecting target spans that are more similar to spans in the datastore) improves detection performance. Notably, across all four domains, the lowest AUROC is 98.5%, suggesting that changes in  $\alpha$  do not cause a substantial performance drop that would change the ranking of detectors. See Appendix C for consistent trends in all generators.

**Impact of Datastore Size.** In our evaluation, ExaGPT uses the train split as the datastore from which it retrieves the top- $k$  most similar spans for each span in a target text. To study the robustness of ExaGPT to datastore size, we analyze how detection performance varies as the datastore size changes. The train split contains 2,000 pairs of human-written and LLM-generated texts. We randomly sample {500, 1,000, 1,500, 2,000} pairs

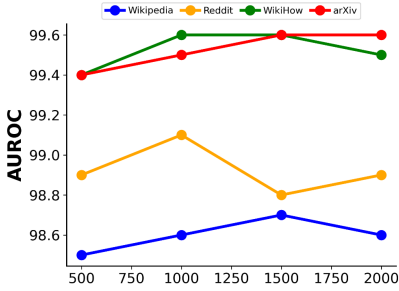


Figure 6: Impact of the datastore size on the detection performance of ExaGPT, using ChatGPT as a generator.

from the train split as datastores of different sizes.

Figure 6 shows the relationship between datastore size and the detection performance of ExaGPT across four domains using ChatGPT as the generator. Overall, ExaGPT remains robust to datastore size, exhibiting only minor performance degradation. Interestingly, ExaGPT with a datastore of 500 pairs performs comparably to using the full 2,000 pairs in terms of AUROC. See Appendix C for consistent trends in all generators.

**Unknown Domain or Generator.** While our primary goal is to improve interpretability in detection, we also perform cross-domain and cross-generator experiments to examine how ExaGPT can be leveraged in more realistic settings. Table 3 reports the cross-domain results. “ALL” denotes the setting where the datastore is constructed from all domains, using equal random samples to match the total size of the single-domain setting. Table 4 reports the cross-generator results. Similarly, “ALL” uses a datastore built from all generators, randomly and equally sampled to match the single-generator size.

We observe that when the domain or generator in our datastore is different from the target model or generator, detection performance is reduced. However, using multiple sources in the datastore is often more realistic in practice, and our results show that in the “ALL” setting, ExaGPT maintains consistently high performance across domains and generators. The tables further offer insights into which domains or generators are effective to consider in the datastore for generalization. For instance, including Reddit mitigates AUROC drops across domains, whereas including ChatGPT helps maintain high detection performance against GPT-4.

**Paraphrased Text.** We also investigate the robustness of ExaGPT against paraphrased text. Following Krishna et al. (2023), we utilize DIPPER,

		Test				
		Wikipedia	Reddit	WikiHow	arXiv	Average
Train	Wikipedia	<b>98.3 / 87.3</b>	91.7 / 68.2	54.1 / 53.3	89.3 / 60.5	83.4 / <b>67.3</b>
	Reddit	90.1 / 60.7	<b>99.3 / 91.1</b>	74.6 / 50.6	93.0 / 63.9	<b>89.3</b> / 66.6
	WikiHow	66.2 / 50.6	76.9 / 60.4	<b>98.8 / 92.2</b>	64.7 / 51.6	76.7 / 63.7
	arXiv	73.7 / 50.4	86.3 / 56.2	57.0 / 51.5	<b>99.7 / 98.7</b>	79.2 / 64.2
	ALL	<u>94.3 / 80.7</u>	<u>96.7 / 83.5</u>	<u>92.9 / 73.4</u>	<u>99.5 / 96.7</u>	<u>95.9 / 83.6</u>

Table 3: Cross-domain detection with GPT-4 as the generator. The scores are AUROC / Acc@FPR=1%.

		Test			
		ChatGPT	GPT-4	Dolly	Average
Train	ChatGPT	<b>99.6 / 95.8</b>	98.2 / 84.5	63.2 / 50.3	87.0 / <b>76.9</b>
	GPT-4	94.6 / 66.6	<b>99.7 / 98.7</b>	61.8 / 51.5	85.4 / 72.3
	Dolly	93.0 / 69.9	89.9 / 65.5	<b>85.2 / 67.3</b>	<b>89.4</b> / 67.6
	ALL	<u>98.9 / 91.4</u>	<u>99.3 / 95.6</u>	<u>76.4 / 52.9</u>	<u>91.5 / 80.0</u>

Table 4: Cross-generator detection with arXiv as the domain. The scores are AUROC / Acc@FPR=1%.

an 11B document-level paraphraser, to rewrite machine-generated text. Table 5 reports results on all domains using ChatGPT as the generator. Among strong interpretable detectors, LR-GLTR was the runner-up in our original in-domain evaluation (§3.2). Even under paraphrasing, ExaGPT maintains moderately high detection performance and consistently outperforms LR-GLTR.

**Comparison with state-of-the-art detectors.** To understand the trade-off between interpretability and detection performance, we also compare ExaGPT with state-of-the-art non-interpretable detectors. Table 6 reports results for Binoculars and Fast-DetectGPT, which are strong metrics-based detectors. The evaluation is conducted on ChatGPT-generated text across all domains. Notably, despite being interpretable, ExaGPT achieves detection performance on par with or even better than these state-of-the-art non-interpretable detectors.

**Inference Cost.** We evaluate the inference efficiency of ExaGPT, focusing on the cost of the  $k$ -NN search. Reducing the datastore from 2,000 to 500 pairs substantially decreases memory usage and inference latency, with almost no loss in detection performance. Using  $k$ -NN approximation further reduces the cost while still maintaining competitive performance. These results indicate that ExaGPT can be deployed efficiently even under limited computational budgets. Detailed results and tables are presented in Appendix C.

Detector	Wikipedia	Reddit	WikiHow	arXiv	Average
LR-GLTR	89.4 / 60.2	97.0 / 76.8	89.5 / 58.0	<b>99.6 / 96.7</b>	93.9 / 72.9
ExaGPT	<b>98.0 / 86.5</b>	<b>97.2 / 69.4</b>	<b>91.1 / 73.8</b>	97.7 / 76.4	<b>96.0 / 76.5</b>

Table 5: Performance on paraphrased text. The scores are AUROC / Acc@FPR=1%.

Detector	Wikipedia	Reddit	WikiHow	arXiv	Average
Binoculars	83.4 / 49.4	72.4 / 55.0	80.9 / 65.0	84.1 / 58.2	80.2 / 56.9
F.GPT	<b>99.6 / 98.1</b>	<b>99.4 / 94.7</b>	<b>95.8 / 85.4</b>	<b>100.0 / 98.1</b>	<b>98.7 / 94.1</b>
ExaGPT	<u>98.6 / 92.3</u>	<u>98.9 / 86.6</u>	<b>99.5 / 96.0</b>	<u>99.6 / 95.8</u>	<b>99.2 / 92.7</b>

Table 6: Performance Comparison with state-of-the-art non-interpretable detectors. The scores are AUROC / Acc@FPR=1%. F.GPT: Fast-DetectGPT.

## 5 Related Work

**LLM-Generated Text Detection.** Prior studies have presented various types of detection algorithms for LLM-generated text, which can be broadly grouped into three categories: *text watermarking*, *metrics-based*, and *supervised classifiers*. Text watermarking modifies the decoding process so that secretly selected tokens appear more frequently, enabling detection by checking their ratio in the output (Kirchenbauer et al., 2023). Metrics-based methods measure the probabilistic discrepancy of a text with the model’s predicted distribution, using signals such as token log probabilities (Gehrmann et al., 2019), token ranks (Solaiman et al., 2019; Su et al., 2023b), entropy (Lavergne et al., 2008), perplexity (Beresneva, 2016; Hans et al., 2024), and negative probability curvature (Mitchell et al., 2023; Bao et al., 2024). Supervised classifiers are models specifically fine-tuned to discern human-written and LLM-generated texts with labels. They vary from probabilistic (Ippolito et al., 2020; Crothers et al., 2023) to neural methods (Uchendu et al., 2020; Rodriguez et al., 2022; Guo et al., 2023).

**Interpretability of the Detection Results.** To minimize the undesired consequences of detection, there is need to develop an LLM detector that provides interpretable evidence for the decision. However, most detectors output only a binary label, and only a few studies aim to provide interpretable evidence. GLTR (Gehrmann et al., 2019) highlights tokens with high model likelihood. Other studies apply explainable machine learning techniques such as LIME and SHAP to supervised classifiers (Mitrović et al., 2023; Wang et al., 2024; Ribeiro et al., 2016; Lundberg and Lee, 2017). DNA-GPT

(Yang et al., 2023) compares  $n$ -gram overlaps between the target text and LLM-generated continuations, providing actual LLM texts with overlaps as evidence.

Unlike prior interpretable detectors, our ExaGPT is grounded by the human decision-making process (Maurer et al., 2006; Barrón-Cedeño et al., 2013) of verifying the origin of a text and can provide more interpretable evidence, as explained in the previous sections.

**Example Retrieval for Interpretability.** Beyond LLM text detection, presenting retrieved examples has been widely used to improve interpretability across NLP tasks, from text generation (Khandelwal et al., 2020) to sequential text classification (Wiseman and Stratos, 2019; Jurafsky et al., 2020; Kaneko et al., 2022). In these approaches, predictions are typically obtained by interpolating the base model’s output distribution with a distribution derived from retrieved nearest-neighbor examples.

Our work has a similar direction of using retrieved similar examples for better interpretability with prior studies in other NLP tasks. In LLM text detection, it is critical to segment the target text into  $n$ -gram spans with individually assigned labels (Cheng et al., 2025). ExaGPT therefore retrieves similar examples for each span and optimizes the final segmentation using dynamic programming.

## 6 Conclusion

We introduced ExaGPT, an interpretable human vs. machine detection approach grounded in the human decision-making process of verifying the origin of a text. In particular, ExaGPT classifies a text by examining whether it shares more verbatim and semantically similar spans with human-written vs. with LLM-generated texts from an available datastore. As evidence of the detection, ExaGPT offers similar span examples for each span in the text. The human evaluation and further analysis show that providing similar span examples allows users to judge the correctness of the detection more effectively than prior interpretable detectors. Moreover, extensive experiments in various domains and generators revealed that ExaGPT has shown notably superior detection performance compared to previous strong detectors, even at a false positive rate of 1%. These results indicate that ExaGPT is a detector with both high interpretability in its decision and high detection performance.

## 7 Limitations

**Bias in the Human Judgments.** Our human evaluation involved four participants with NLP backgrounds, which may limit generalizability to typical users. We selected the evaluators to ensure high-quality, consistent feedback for this initial and information-intensive task, requiring annotators to carefully read and assess the validity of the detection based on its evidence for every single sample. Although ExaGPT is designed to be intuitive even for non-experts by providing concrete span-level textual evidence rather than probabilistic or attribution scores, we leave its validation on broader user populations for future work.

## 8 Ethics and Broader Impact

**Human Subject Considerations.** In our study, human subjects are engaged in identifying the correctness of the detection based on evidence. All annotators provided informed consent, were fully aware of the study’s objectives, and had the right to withdraw at any time.

**Transparency and Reproducibility.** To promote open research, we release our code and data to the public, including all human annotations.

## References

Guangsheng Bao, Yanbin Zhao, Zhiyang Teng, Linyi Yang, and Yue Zhang. 2024. [Fast-detectgpt: Efficient zero-shot detection of machine-generated text via conditional probability curvature](#). *Preprint*, arXiv:2310.05130.

Alberto Barrón-Cedeño, Marta Vila, M. Antònia Martí, and Paolo Rosso. 2013. [Plagiarism meets paraphrasing: Insights for the next generation in automatic plagiarism detection](#). *Computational Linguistics*, 39(4):917–947.

Daria Beresneva. 2016. Computer-generated text detection using machine learning: A systematic review. In *21st International Conference on Applications of Natural Language to Information Systems, NLDB*, pages 421–426. Springer.

Bloomberg. 2024. [Ai detectors falsely accuse students of cheating—with big consequences](#). Accessed on 2024-10-20.

Zihao Cheng, Li Zhou, Feng Jiang, Benyou Wang, and Haizhou Li. 2025. Beyond binary: Towards fine-grained LLM-generated text detection via role recognition and involvement measurement. In *THE WEB CONFERENCE 2025*.

Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. [Free Dolly: Introducing the World’s First Truly Open Instruction-Tuned LLM](#). Accessed: 2024-7-12.

Evan Crothers, Nathalie Japkowicz, and Herna Viktor. 2023. [Machine generated text: A comprehensive survey of threat models and detection methods](#). *Preprint*, arXiv:2210.07321.

Liam Dugan, Alyssa Hwang, Filip Trhlík, Andrew Zhu, Josh Magnus Ludan, Hainiu Xu, Daphne Ippolito, and Chris Callison-Burch. 2024. [RAID: A shared benchmark for robust evaluation of machine-generated text detectors](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12463–12492, Bangkok, Thailand. Association for Computational Linguistics.

Sebastian Gehrmann, Hendrik Strobelt, and Alexander M. Rush. 2019. [Gltr: Statistical detection and visualization of generated text](#). *Preprint*, arXiv:1906.04043.

Gizmodo. 2024. [AI Detectors Get It Wrong. Writers Are Being Fired Anyway](#). Accessed on 2024-07-12.

Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jiran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. [How close is chatgpt to human experts? comparison corpus, evaluation, and detection](#). *Preprint*, arXiv:2301.07597.

Abhimanyu Hans, Avi Schwarzschild, Valeriia Cherepanova, Hamid Kazemi, Aniruddha Saha, Micah Goldblum, Jonas Geiping, and Tom Goldstein. 2024. [Spotting llms with binoculars: Zero-shot detection of machine-generated text](#). *Preprint*, arXiv:2401.12070.

Daphne Ippolito, Daniel Duckworth, Chris Callison-Burch, and Douglas Eck. 2020. [Automatic detection of generated text is easiest when humans are fooled](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1808–1822, Online. Association for Computational Linguistics.

Jiazhou Ji, Ruizhe Li, Shujun Li, Jie Guo, Weidong Qiu, Zheng Huang, Chiyu Chen, Xiaoyu Jiang, and Xinru Lu. 2024. [Detecting machine-generated texts: Not just "ai vs humans" and explainability is complicated](#). *Preprint*, arXiv:2406.18259.

Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. [Billion-scale similarity search with gpus](#). *Preprint*, arXiv:1702.08734.

Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault. 2020. Proceedings of the 58th annual meeting of the association for computational linguistics. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.

645	Masahiro Kaneko, Sho Takase, Ayana Niwa, and Naoaki Okazaki. 2022. <a href="#">Interpretability for language learners using example-based grammatical error correction</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7176–7187, Dublin, Ireland. Association for Computational Linguistics.	699
646		700
647		701
648		702
649		703
650		704
651		705
652	Urvashi Khandelwal, Angela Fan, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Nearest neighbor machine translation. <i>arXiv preprint arXiv:2010.00710</i> .	706
653		707
654		708
655		709
656		710
657	John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2023. <a href="#">A Watermark for Large Language Models</a> . <i>Preprint</i> , arXiv:2301.10226.	711
658		712
659		713
660	Ryuto Koike, Masahiro Kaneko, and Naoaki Okazaki. 2024. <a href="#">OUTFOX: LLM-Generated Essay Detection Through In-Context Learning with Adversarially Generated Examples</a> . In <i>Proceedings of the 38th AAAI Conference on Artificial Intelligence</i> , Vancouver, Canada.	714
661		715
662		716
663		717
664		718
665		719
666	Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Wieting, and Mohit Iyyer. 2023. <a href="#">Paraphrasing evades detectors of ai-generated text, but retrieval is an effective defense</a> . <i>Preprint</i> , arXiv:2303.13408.	720
667		721
668		722
669		723
670	Thomas Lavergne, Tanguy Urvoy, and François Yvon. 2008. <a href="#">Detecting Fake Content with Relative Entropy Scoring</a> . In <i>Proceedings of the ECAI'08 Workshop on Uncovering Plagiarism, Authorship and Social Software Misuse</i> , CEUR Workshop Proceedings.	724
671		725
672		726
673		727
674		728
675	Yinhan Liu. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> , 364.	729
676		730
677		731
678		732
679	Scott Lundberg and Su-In Lee. 2017. <a href="#">A unified approach to interpreting model predictions</a> . <i>Preprint</i> , arXiv:1705.07874.	733
680		734
681	Hermann Maurer, Frank Kappe, and Bilal Zaka. 2006. Plagiarism – a survey. <i>Journal of Universal Computer Science</i> , 12(8):1050–1084.	735
682		736
683		737
684	Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. 2023. <a href="#">DetectGPT: Zero-Shot Machine-Generated Text Detection using Probability Curvature</a> . <i>Preprint</i> , arXiv:2301.11305.	738
685		739
686		740
687		741
688		742
689	Sandra Mitrović, Davide Andreoletti, and Omran Ayoub. 2023. <a href="#">Chatgpt or human? detect and explain explaining decisions of machine learning model for detecting short chatgpt-generated text</a> . <i>Preprint</i> , arXiv:2301.13852.	743
690		744
691		745
692		746
693		747
694	OpenAI. 2023a. <a href="#">How can educators respond to students presenting ai-generated content as their own?</a> Accessed: 2024-6-10.	748
695		749
696		750
697	OpenAI. 2023b. <a href="#">Introducing ChatGPT</a> . Accessed on 2024-03-10.	751
698		752
		753
		754
		755
		756
	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. <a href="#">"why should i trust you?": Explaining the predictions of any classifier</a> . <i>Preprint</i> , arXiv:1602.04938.	757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800

757	Sam Wiseman and Karl Stratos. 2019. <a href="#">Label-agnostic sequence labeling by copying nearest neighbors</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5363–5369, Florence, Italy. Association for Computational Linguistics.	805
758		806
759		
760		
761		
762		
763	Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan, Derek F. Wong, and Lidia S. Chao. 2023. <a href="#">A survey on llm-generated text detection: Necessity, methods, and future directions</a> . <i>Preprint</i> , arXiv:2310.14724.	810
764		811
765		812
766		813
767	Xianjun Yang, Wei Cheng, Yue Wu, Linda Petzold, William Yang Wang, and Haifeng Chen. 2023. <a href="#">Dna-gpt: Divergent n-gram analysis for training-free detection of gpt-generated text</a> . <i>Preprint</i> , arXiv:2305.17359.	814
768		815
769		816
770		817
771		818
772	<b>A Detailed Configurations of Baselines</b>	819
773	<b>LR-GLTR.</b> Following the setting of (Wang et al., 2024), we leverage the two categories of GLTR features: (1) the number of tokens in the top- $\{10, 100, 1,000, 1,000+\}$ ranks in the predicted probability distribution of LLMs (four features), and (2) the probability distribution of the word divided by the maximum probability of any word at the same position over 10 bins between 0.0 and 1.0 (ten features).	820
774		821
775		822
776		823
777		824
778		825
779		826
780		827
781		828
782	<b>DNA-GPT.</b> For DNA-GPT, we set the truncation ratio $\gamma$ to 0.7 and 0.5, and the number of re-generations $K$ to 10 and 5 for closed-source and open-source LLMs. We ensured that the temperature is the same as the one used to generate a target text and that the generation prompt is known. These configurations were found to ensure the favorable performance of DNA-GPT in (Yang et al., 2023). We set all other hyperparameters to their default values.	829
783		830
784		831
785		832
786		
787		
788		
789		
790		
791		
792	<b>B Detection Evidence of Baselines</b>	
793	<b>RoBERTa with SHAP.</b> Figure 7 depicts an example of evidence by RoBERTa with SHAP. We visualize the evidence using the SHAP library <sup>9</sup> . Overall, the red parts are spans that contribute to predicting LLM-generated. The blue parts are spans that contribute to predicting human-written. In the evidence, if the prediction value, $f(\text{inputs})$ moves further to the right compared to the base value (the expected value across all data samples), it is more likely to be LLM-generated. When we hover over a colored part, we can also see a score of how much the part contributes to the detection	833
794		834
795		835
796		836
797		
798		
799		
800		
801		
802		
803		
804		
	result. The more a span contributes to the decision, the darker its color.	
	<b>LR-GLTR.</b> Figure 8 displays an example of evidence by LR-GLTR. We leverage a demo app <sup>10</sup> of GLTR, provided by Gehrmann et al. (2019). It highlights tokens in different colors based on their rank of top- $\{10, 100, 1,000, 1,000+\}$ in the predicted token distribution from an LLM. The higher the rank of the token, the more likely an LLM is to generate the token. The green parts are spans that an most likely LLM-generated. The degree decreases in the order of green, yellow, red, and purple. When we hover a cursor on a colored part, we can also see the predicted token distribution of an LLM.	
	<b>DNA-GPT.</b> Figure 9 shows an example of evidence by DNA-GPT. We implemented a demo app of DNA-GPT with the streamlit framework <sup>11</sup> . It shows overlapped $n$ -gram spans between a truncated target text and multiple LLM-generated continuations. The more blue spans, the more likely the text is LLM-generated. For span matching, we follow the original implementation of DNA-GPT <sup>12</sup> where it was achieved by token-level matching based on preprocessing of the lower casing and stemming. We also set $n$ to 8 in order to show a large number of overlapped spans enough to interpret as evidence.	
	<b>C Analysis Details</b>	
	<b>Example of <math>k</math>-NN spans.</b> Table 7 presents examples of long spans ( $n = 19$ ) with high reliability scores for a target span retrieved by ExaGPT.	
	<b>Impact of <math>\alpha</math>.</b> Figure 10 showcases the impact of $\alpha$ on the detection performance of ExaGPT across four domains and three generators. We found similar overall trends of the impact of $\alpha$ in other LLMs, including GPT-4 and Dolly-v2, with the impact in ChatGPT, as explained in §4.	
	<b>Impact of the Datastore Size.</b> Figure 11 showcases the impact of the datastore size on the detection performance of ExaGPT across four domains and three generators. We can observe similar overall trends of the impact of datastore size in other LLMs, including GPT-4 and Dolly-v2, with the impact in ChatGPT as explained in §4.	

<sup>9</sup><https://shap.readthedocs.io/>

<sup>10</sup><http://demo.gltr.io/client/index.html>

<sup>11</sup><https://github.com/streamlit/streamlit>

<sup>12</sup><https://github.com/Xianjun-Yang/DNA-GPT>

Target Span	LLM	published in 1993. The novel tells the story of a young Jewish slave, Hadassah,
<i>k</i> -NN Spans	LLM (0.92)	and was first published in 1936. The book tells the story of three orphaned sisters,
	LLM (0.92)	published in 2012. The novel revolves around the story of a young woman
	LLM (0.90)	and published in 2010. The novel tells the story of Michael Beard, a
	LLM (0.90)	ling of the biblical book, Song of Solomon, and is considered one of the
	LLM (0.90)	man and published in 1963. The book was later adapted into a Disney film of the
	LLM (0.90)	. The film tells the story of a young
	Human (0.89)	the Xanth series. It is the second book of a trilogy beginning with Vale of the
	LLM (0.89)	published in 1959. The novel is set in the Arctic region and follows the story of Dr.
	Human (0.89)	. It is the third novel in the Dahak trilogy, after the de
LLM (0.89)	for his semi-autobiographical novel, “The Watch that Ends the Night”. Born in	

Table 7: Examples of *k*-NN spans for a target span retrieved by ExaGPT. The colored part represents the original label for each span (LLM in blue and Human in red, respectively). In the part of *k*-NN spans, the similarity between the target span and each *k*-NN span is added.

	#Instance	GPU memory (GB)	Latency (sec.)	AUROC
2000 pair	36M	162.2	14.6	99.5
500 pair	9.1M	54.7 (66%↓)	5.81 (60%↓)	99.4
500 pair + IVFPQ	9.1M	20.2 (87%↓)	1.22 (90%↓)	97.8

Table 8: Inference cost analysis of ExaGPT.

**Inference Cost.** In our preliminary analysis, we found that the inference cost of embedding generation and DP-based segmentation was negligible compared to the *k*-NN search, which was the primary bottleneck due to the large datastore containing extensive *n*-gram instances. Therefore, we have measured inference latency and GPU memory usage focusing on the *k*-NN search component. We also have conducted experiments with FAISS using IVFPQ indexes, to reduce resource usage and improve inference speed.

Table 8 provides the results on the WikiHow domain with ChatGPT. Here, #Instance refers to the number of *n*-gram spans in the datastore. While we observe that achieving full performance with a datastore of 2,000 pairs requires considerable computational resources and higher latency, reducing the size to 500 pairs decreases GPU memory usage by 66% and latency by 60% without compromising detection performance. With the additional use of FAISS-based *k*-NN approximation, the requirements are further reduced by 87% in memory and over 90% in latency, while the performance drop is still moderate. These findings highlight the promising practical applicability of ExaGPT.

## D Computational Budget

We run all the experiments with two AMD EPYC 7453 CPUs and four NVIDIA A6000 GPUs. The

total processing time is approximately 25 hours.

878

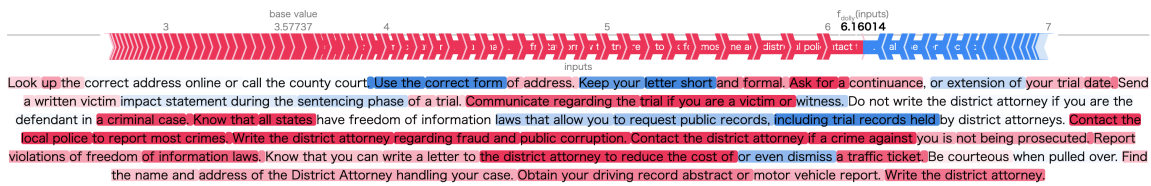
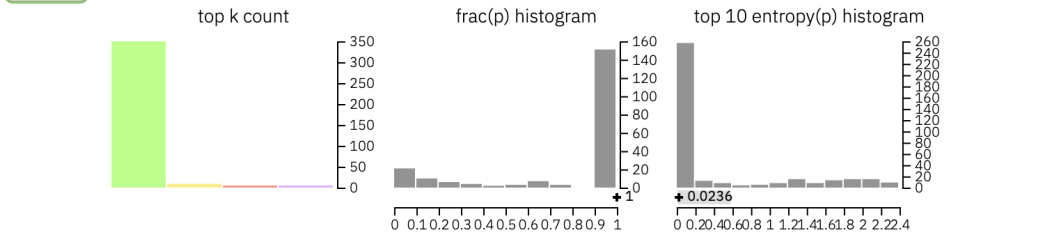


Figure 7: Example of evidence by ROBERTa with SHAP.

or enter a text:

Thomas C. Hinkle was an American politician and lawyer. He was a member of the Republican Party and served as the U.S. Representative for, from 1991 to 1995. He ran for the U.S. Senate from Louisiana in 1996 but lost to Mary Landrieu. Hinkle is a resident of Rowlett, Texas and was a news anchor and reporter for CBS affiliate station KXAS-TV in Fort Worth from 1980 to

analyze



Top K: 10 100 1000 10000

Thomas C. Hinkle was an American politician and lawyer. He was a member of the Republican Party and served as the U.S. Representative for, from 1991 to 1995. He ran for the U.S. Senate from Louisiana in 1996 but lost to Mary Landrieu. Hinkle is a resident of Rowlett, Texas and was a news anchor and reporter for CBS affiliate station KXAS-TV in Fort Worth from 1980 to 1991. discuss Thomas C. Hinkle was an American politician and lawyer. He was a member of the Republican Party and served as the U.S. Representative for, from 1991 to 1995. He ran for the U.S. Senate from Louisiana in 1996 but lost to Mary Landrieu. Hinkle is a resident of Rowlett, Texas and was a news anchor and reporter for CBS affiliate station KXAS-TV in Fort Worth from 1980 to 1991. evidentiary Thomas C. Hinkle was an American politician and served as the U.S. Representative for, from 1991 to 1995. He ran for the U.S. Senate from Louisiana in 1996 but lost to Mary Landrieu. Hinkle is a resident of Rowlett, Texas and was a news anchor and reporter for CBS affiliate station KXAS-TV in Fort Worth from 1980 to 1991. happenings Thor and served as the U.S. Representative for, from 1991 to 1995. He ran for the U.S. Senate from Louisiana in 1996 but lost to Mary Landrieu. Hinkle is a resident of Rowlett, Texas and was a news anchor and reporter for CBS affiliate station KXAS-TV in Fort Worth from 1980 to 1991. top\_k pos: 0 prob: 0.999 frac(p): 1.000

Figure 8: Example of evidence by LR-GLTR.

templat the possibl option of an apprenticeship or on the job train program an apprenticeship or on the job train program is an excel way to gain hand on experi in the automot technician trade while you earn money consid thi option if you prefer a more hand on learn experi 10 contact your guidanc counselor if you are a high school student consid enter the automot technician trade if you are a high school student talk to your guidanc counselor about the vocat train program avall in your area they can provid guidanc on what program are best suit for your goal and interest 11 visit sever websit that offer educ train in the automot technician trade research variou vocat train program onlin to compar offer requir and cost make sure you read review from previou student to see what they have to say about the program 12 select the best program for your need interest and career goal after review your option choos the program that best fit your need interest and career goal 13 chart out your final plan and proceed with confid onc you have select your program and commit to your career plan out your step to achiev your goal stay motiv and confid in your abil to succeed as an automot technician

0. templat the possibl option of an apprenticeship or on the job train program depend on the program apprenticeship and on the job train program a great way to gain practic experi and learn from experienc profession while earn an incom consid if thi is the right path for you 10 contact your guidanc counselor if you are a high school student consid enter the automot technician trade high school often offer vocat train program in automot technolog speak with your guidanc counselor to learn more about these program and if they are avall to you 11 visit sever websit that offer educ train in the automot technician trade research onlin for differ vocat train program and read through their program curriculum determin which program will best suit your desir career path 12 select the best program for your need interest and career goal choos the program that align with your career aspir and offer the level of train and certifi you desir 13 chart out your final plan and proceed with confid now that you ve research and weigh your option it is time to start your train program stick to your plan and be confid in your decis to become an automot technician with hard work and dedic you can have a success career in the field of automot technolog

1. templat the possibl option of an apprenticeship or on the job train program apprenticeship and on the job train program provid a more hand on approach to learn the automot technician trade consid thi option if you prefer a more practic and immers learn environ 10 contact your guidanc counselor if you are a high school student consid enter the automot technician trade your guidanc counselor can provid you with valuabl resourc and inform about vocat train program and opportun 11 visit sever websit that offer educ train in the automot technician trade research differ vocat train program and school that offer automot technolog cours look for program that have a good reput and offer hand on train opportun 12 select the best program for your need interest and career goal choos a program that fit with your learn style schedul and goal look for program that provid a well round educ in differ area of automot technolog 13 chart out your final plan and proceed with confid onc you have select a program make a plan to achiev your goal and proceed with confid stay commit to your train and take advantag of any opportun to gain hand on experi and network with other profession in the field conclus become an automot technician requir dedic skill and a willing to learn research the field evaly your strength and weak and determin the best type of vocat train program for your need consid your goal and determin if you want to special in a particular area of the trade or pursu a

Figure 9: Example of evidence by DNA-GPT.

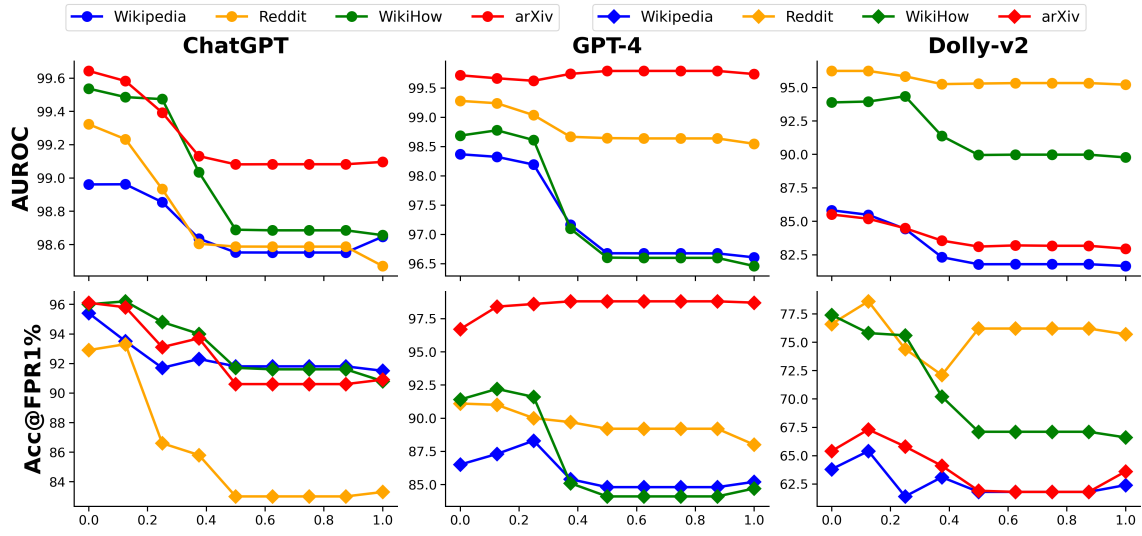


Figure 10: Impact of  $\alpha$  on the detection performance of ExaGPT, including the AUROC and the accuracy at 1% FPR, across four domains and three generators.

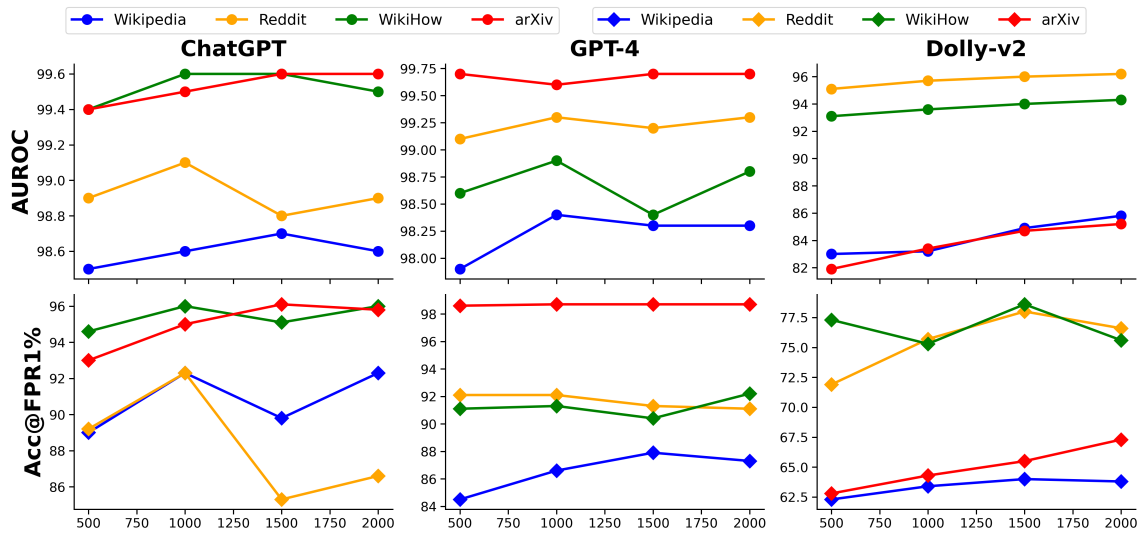


Figure 11: Impact of the datastore size on the detection performance of ExaGPT, including the AUROC and the accuracy at 1% FPR, across four domains and three generators.